# DØ  PPDG Year 3 plans

**Table of Contents**

# DØ Team

## 1 Introduction

The DØ-PPDG project is the core of the DØGrid (now SAMGrid). It's continuing mission to provide truly distributed (in the global sense) computing for the Run II experiments. We leverage SAM – a mature, advanced data management system, and add services for job and information management, while collaborating with computer scientists, most notably, University of Wisconsin. Our principal deliverable in SAMGrid is secure, reliable execution of structured, globally distributed jobs, with sufficient provision of monitoring, both at run time and historically.

During the first two years of our participation in the collaboration, we have delivered the first version of JIM (the Job and Information Management of SAMGrid), where we integrated the power of Condor-G for Grid job brokering with that of SAM for global data movement. Our global scheduler dispatches jobs based on the amount of data cached at participating sites. Primary job type is analysis, with rudimentary support for Monte-Carlo and reconstruction. In addition, web-based monitoring of the whole system and individual jobs was provided. Presently, JIM V1 is in the process of deployment.

## 2 PPDG Year 3 Plans

During the third year of the collaboration, we plan to stabilize JIM v1, expose it to high operational load. We will define and deliver Version 2 of JIM, as well as prototype an operational support model for the (remote) SAMGrid installations.

Planned features for JIM (6-9 mos):

- Full support for MC and reconstruction jobs for both experiments

- Full support for logging of all "interesting" events, filtering and archiving

- Next generation VOMS, possibly imported from Europe

The following developments will have to take place in JIM:

- Resource description for MC environments has to be frozen and published, using the JIM-adopted Condor ClassAd framework. Presently, a lot of ideas have been circulated and a prototype implementation is available. The same applies to reconstruction.

- The DMZ between the Grid and the Fabric has to be clearly defined. One of the novelties of JIM's design is the use of the XML database on the border between the Grid and the Fabric. We will need to expose this topic to a broader Grid community so as to stimulate discussions and possibly standardize the interface between the Fabric and the Grid.

- The XML-based logger has to be fully developed. Unlike the Logging and Book-keeping service from EDG WP1, our logging service is an important concept which will underlie historical data mining for both data transfers and jobs, and thus provide accounting at various levels.

- Evaluation of Web services must also be complete by the end of year 3. While any particular grid-like system like SAM can do without web services, it is becoming increasingly clear that it is impossible to combine diversity in grid solutions with grid interoperability without a common language to describe services, and WSDL is the *de facto* such leader. Thus, unless the Grid community in general and Run II experiments in particular want to entrench into middleware consisting of proprietary (GTK 2) implementations of obscure, often criticized protocols (GRAM, GridFTP), we must arm ourselves with proper generalizations and describe our system in a language like WSDL.

The above work will keep the two PPDG FTE's (with helper students) well occupied. In addition, more work will be required to SAM. However, we envision more intense evolution of this data handling system towards the Grid.

Thus, the evolution of SAM needs to be linked more tightly with PPDG activities and we will make according proposals for years 4 and 5 below.

Currently, SAM has undergone a significant change whereby its core data replication mechanism has been disentangled from local cache management, thus allowing for virtually arbitrary storage system to serve as a source of consumers' data and/or a "local" cache. Independently of PPDG progress during year 3, and simultaneously with it, we expect the process of modularization to continue, so as to generate a *bi-directional* idea exchange between storage resource management frameworks (SRM etc) and core SAM's data replicator. This will invariably have a bearing on computational economics of scheduling of data-intensive jobs.

## 2.1 Action Items from '03 Questionnaires

We continue to work closely with the Condor team to develop and test  the three-tier condor architecture, and test the matchmaking features provided for JIM by the Condor team.  We ask that these efforts continue to be high priority for the Condor team. An installation of a JIM Gateway is in progress and this will facilitate continued cooperation between the JIM and Condor teams, and further  understanding, and testing of the relevant software.

## 2.2 Project JIM

We are now in the process of deploying the first production version of the JIM product for DØ and CDF.  Success of the software is measured by first providing the experiments with functionality  currently  used by the experiments for accomplishing their computing tasks.  The next step is  enabling the automated job management features which provide transparent access to new resources, and reduced manpower requirements to manage these activities.

### 2.2.1 Deliverables

The following list of deliverables are items the team, DØ and CDF feel are needed to make the system fully functional for  the Run 2 experiments. It is not necessarily in prioritized order, and not all the items in the list will be achievable during year three.

1. Deployment of JIM
   a. Sites (execution, submission, monitoring)
      i. Initial DØ sites (GridKa, Wupprtal, Lancaster, IC, RAL, Lyon, NIKHEF, U. Michigan, U Wisconsin)
      ii. CDF sites (GridKa, Glasgow, Oxford, Scotgrid, Trieste, Toronto, Kyungpook, TexasTech, UCSD)
      iii. Extended DØ site list (UTA, Manchester, Prague, potentially any existing SAM station)
      iv. Extended DØ site list for submission (potentially any DØ site)
   b. Support for
      i. Monte Carlo production
      ii. Production processing (data reconstruction)
      iii. User analysis
   c. Provided functionality
      i. Phase 1: Same as existing functionality with manual decisions about site execution
      ii. Phase 2: Enable true brokering of jobs. Means work with condor match making, ranking jobs, resource advertising.
      iii. Phase 3: Structured jobs split intelligently and managed in parallel

2. Issues with input sandbox

    a. Size requirements for gateway node need to be understood

    b. Transfer to worker execution node when the job runs.

3. Return output files from jobs

    a. Transient sam storage location

    b. Alternative techniques, e.g via the output sandbox ( authentication proxy issues)

4. Improve installation framework

    a. Reduce installation time by clarifying the configuration process (interaction with the installer) and developing more intelligent guesses for the defaults

    b. The planned framework is not ready and has caused installation delays, need to complete this.

    c. Packaging and code distribution with ups/upd needs to be evaluated

    d. Improve upgrade procedure

5. Additional troubleshooting tools + test of "sanity" procedure (e.g. after installation)

6. Hardening of the system

    a. Production experience

    b. Extensive load testing with test harness application

    c. upgrading to new/more robust versions of globus and condor

7. Continuing work with Condor team

    a. Continue to improve 3-tier architecture

        i. Improved authentication to provide true 3-tier functionality

        ii. Return of standard out and standard error

    b. Other debugging as problems are observed

8. Configuration management

    a. Have XML framework that is now used

    b. Can easily advertise configuration through web interface

9. Workflow management

    a. Job manager interface for mcrunjob

    b. Cooperation with CMS, FNAL-CD, UK players in mcrunjob (a.k.a. Shahkar project)

    c. Investigate redundancy/fail-over models

10. SRM used in general cache and storage interface strategy for SAM

11. VO Management, in conjunction with one or more of following

    a. Existing JIM VOMS (written by Gabriele Garzoglio)

    b. EDG VOMS project

    c. FNAL/USCMS/SDSS VOX project.

12. logging service

    a. XML based

    b. Support for both reliable and unreliable message delivery

    c. Hierarchical model

13. Monitoring service

    a. Current MDS implementation is probably not adequate (explore alternatives)

    b. Continue to develop Web based interface

    c. (CDF requests access to remote node through unix like commands, ls, cat, etc. )

14. JDL: Job definition language for structured jobs

    a. What is a job?

    b. How is it parsed?

    c. How is it related to the condor?

15. Evaluation of Web services

    a. Condor will move to WS

    b. Globus will move to WS

    c. Develop WS interfaces for SAMGrid

        i. Interface to EDG through WSDL/SOAP adapter for SAM

        ii. (D0 Europe is very interested and may supply manpower)

        iii. Additional WS as needed

16. Gathering statistics on relevant metrics: understand what parameters are minimized/maximized using the (SAM-)Grid paradigm to computation

17. Distributed Replica Catalogues

18. Glue Schema

    a. Continued work on common areas

    b. Implementation where relevant.


**2.2.2 Milestones**

7/1/2003:  Finish  installation of JIM V1.0 at initial sites; begin load testing; exercise MC production;  Resolve issues with user input sandbox management; Evaluation of  VO management options; Begin building load testing proceedures

10/1/2003: JIM V1.1 release; Installation framework ready; Use of VDT for release; Fully functional configuration management;  Improved monitoring; Full three-tier condor implementation; stdout and stderr returned through condor pool; Second round of installs;  Additional troubleshooting toolkit; Initial logging service features. Establish an operational support model. Integration of new VO management tools; Initial phase WS interface for EDG/LCG to SAM RC (not fully approved).


1/1/2003: JIM V1.2 Full logging service features ready;  JDL for structured jobs; Additional scheduling functionality; augmented resource advertisement; Begin further evaluation of WS;  SAM transient file transfer features ready; Improved monitoring system;  Additional workflow management functionality; SRM implementation (at some level, probably for access to dCache, possibly for cache management strategy); Third round of deployment; Updated support features; Improved reliability.

4/1/2003: JIM V1.3; Move toward Globus and Condor WS implementations (pending their availability status); ; Full accounting statistics available; Contingency for previous milestones.;  Improved support and reliability.


7/1/1003: JIM V 1.4; Begin exploring distributed replica catalog; More work on Web services;

### 2.2.3 Participants

Igor Terekhov -  Technical Team Lead

Lee Lueking – DØ Team Lead

Gabriele Garzoglio

Andrew Baranovski

Vijay Murthi – CS Masters student through contract with U. Texas Arlington

Parag Mhashilkar – CS Masters student through contract with U. Texas Arlington

### 2.2.4 Dependencies

Condor matchmaking, condor three-tier architecture. Globus TK,. VDT.

### 2.2.5  Issues and Concerns

The JIM team is faced with many large and challenging tasks in  several CS areas, and we are seriously looking for additional collaboration with other physics teams. We need to work together to get these things done. Some areas where common solutions might be explored include:

- Drawing the boundary between the Grid and Fabric

- Sophisticated logging service using standard technologies like XML-SOAP

- Framework for user Job Definition Language. Note that this is different from a single JDL.

- Framework for Configuring site and its resources (for publishing or discovery). Note that this is different from converging on a fixed schema.

- Intra-cluster workload management and its interface to the Grid

# 3 PPDG Year 4 and 5

Generally, future work will proceed along two avenues: 1) major feature enrichments (robustness, ease of use, better system control), and 2) co-existance and  interoperability with other grids.

As far as JIM is concerned, we propose:

- Manage *structured* jobs at the global level, including decomposition and recombination of job's fragments according to the job details known to the scheduler. Presently, this cumbersome job is done manually by e.g. MC request executor. Ultimately, such decomposition should take place dynamically, i.e. as more resources become available. Some efforts in this area are undertaken at the local cluster level, moving it to the grid level will present new challenges.

- Tolerance to entire grid cluster/site failure, with re-routing of the corresponding job fragments to new resources.

- The above two may collectively be called workflow management, but one should not confuse intra-cluster management with that at the grid level.

For inter-operability, we intend to pursue studies of the Glue schema and Job description languages. Both need to be looked at from a broader prospective that defining e.g. standard UML diagram or a standard language. We are more inclined to converging on a framework that allows to process data representing these entities. We are also working with various other Grid projects to define common interfaces, generally as web services. One example is the work started with EDG WP2 to explore such an interface.

As far as SAM per se is concerned, in addition to the previously mentioned enhancements, we will need the following:Entanglement of all meta-data into a single Oracle database has to change. Some data items are critical for the entire collaboration and must be conceptually centralized. Other items are of highly transient nature and hinder deployment of large distributed caches at distributed sites. We will likely need splitting of the schema, decentralization that the catalog person will be confident bookkeeping service about the bigger (of SAM towards HV, as well as for the overall success of D0-PPDG such as replica catalogue) with externally developed, more standard ones.

# 4 Background

## 4.1 Summary Table from PPDG proposal

This is the summary table from the PPDG proposal. The full proposal is available at
http://lbnl2.ppdg.net/docs/scidac01_ppdg_public.doc

We have added a "+" in the Yr3 column for D0 team activities in these areas

| Project Activity | Experiments | Yr1 | Yr2 | Yr3 |
|---|---|---|---|---|
| CS-1 Job Description Language – definition of job processing requirements and policies,  file placement & replication in distributed system. | | | | |
| P1-1 Job Description Formal Language | D0, CMS | X | | + |
| P1-2 Deployment of Job and Production Computing Control | CMS | X | | + |
| P1-3 Deployment of Job and Production Computing Control | ATLAS, BaBar, STAR | | X | |
| P1-4 Extensions to support object collections, event level access etc. | All | | | X |
| CS-2 Job Scheduling and Management   - job processing, data placement, resources discover and optimization over the Grid | | | | |
| P2-1 Pre-production work on distributed job management and job placement optimization techniques | BaBar, CMS, D0 | X | | + |
| P2-2 Remote job submission and management of production computing activities | ATLAS, CMS, STAR, JLab | | X | + |
| P2-3 Production tests of network resource discovery and scheduling | BaBar | | X | |
| P2-4 Distributed data management and enhanced resource discovery and optimization | ATLAS, BaBar | | | X |
| P2-5 Support for object collections and event level data access.  Enhanced data re-clustering and re-streaming services | CMS, D0 | | | X |
| CS-3 Monitoring and Status Reporting | | | | |
| P3-1 Monitoring and status reporting for initial production deployment | ATLAS | X | | |
| P3-2 Monitoring and status reporting – including resource availability, quotas, priorities, cost estimation etc | CMS, D0, JLab | X | X | + |
| P3-3 Fully integrated monitoring and availability of information to job control and management. | All | | X | X |
| CS-4 Storage resource management | | | | |

| | | | | |
|---|---|---|---|---|
| P4-1  HRM extensions and integration for local storage system. | ATLAS, JLab, STAR | X | | + |
| P4-2 HRM integration with HPSS, Enstore, Castor using GDMP | CMS | X | | + |
| P4-2  Storage resource discovery and scheduling | BaBar, CMS | | X | |
| P4-3 Enhanced resource discovery and scheduling | All | | | X |
| CS-5 Reliable replica management services | | | | |
| P5-1 Deploy Globus Replica Catalog services in production | BaBar, JLab | X | | |
| P5-2 Distributed file and replica catalogs between a few sites | ATLAS, CMS, STAR, JLab | X | | |
| P5-3  Enhanced replication services including cache management | ATLAS, CMS | | X | |
| CS-6 File transfer services | | | | |
| P6-1 Reliable file transfer | ATLAS , BaBar, CMS, STAR, JLab | X | | |
| P6-2 Enhanced data transfer and replication services | ATLAS, BaBar, CMS, STAR, JLab | | X | |
| CS-7 Collect and document current experiment practices and potential generalizations | All | X | X | X |

## 4.2 Evaluation and Continuing Support

Also from the proposal. Please could you comment on these items from your teams perspective:

"During the last phase of PPDG we will: (not sure how to answer these questions)

- Transition all technologies that remain in service to a sustaining support mode. Support agreements will be explored with commercial vendors as well as with the collaborating CS groups, proposed GridPPs team, and the experiments.
- Document the pros and cons of the  approaches used by  this project and propose follow-on projects after it's completion.
- Attempt to further generalize the common approach to multiple physics applications into recommendations for common approaches to multiple science areas.
- Explore with software vendors possible commercialization of appropriate components of PPDG software. "